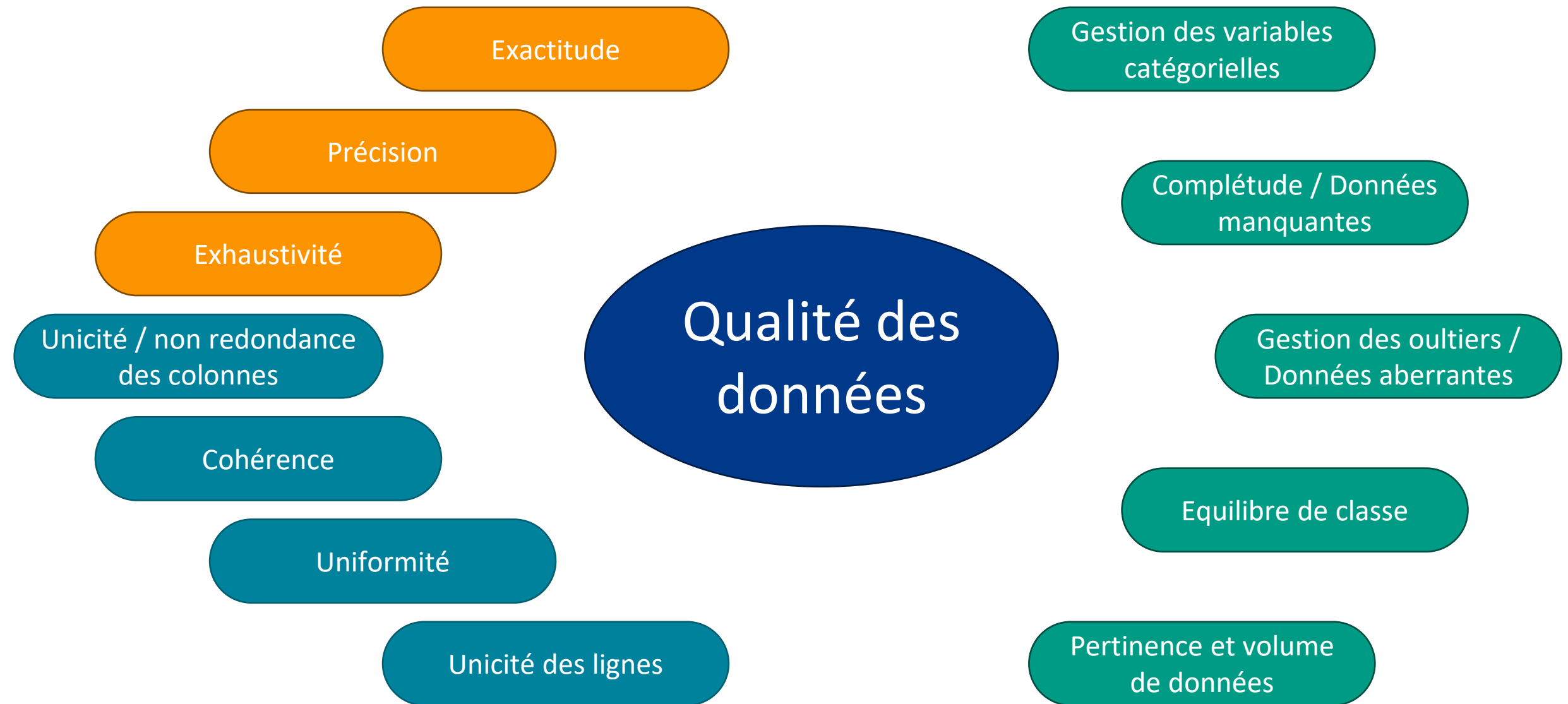




02

## Atelier 2

# Atelier 2 : La qualité des données, un problème multiple



# Atelier 2 : Normalisation

Dans un jeu de données tabulaires, on retrouve fréquemment des données catégorielles. On distingue :

Les variables nominales (sans ordre) :

- **Type de véhicule** : Voiture, Camion, Moto, Vélo
- **Nationalité** : FR, EN, ES
- **Etat civil** : Célibataire, Marié, Divorcé, Veuf

## One-hot Encoding

Nationalité	Nat-FR	Nat-EN	Nat-ES
FR	1	0	0
FR	1	0	0
EN	0	1	0
ES	0	0	1
FR	1	0	0
EN	0	1	0

Les variables ordinales (avec ordre) :

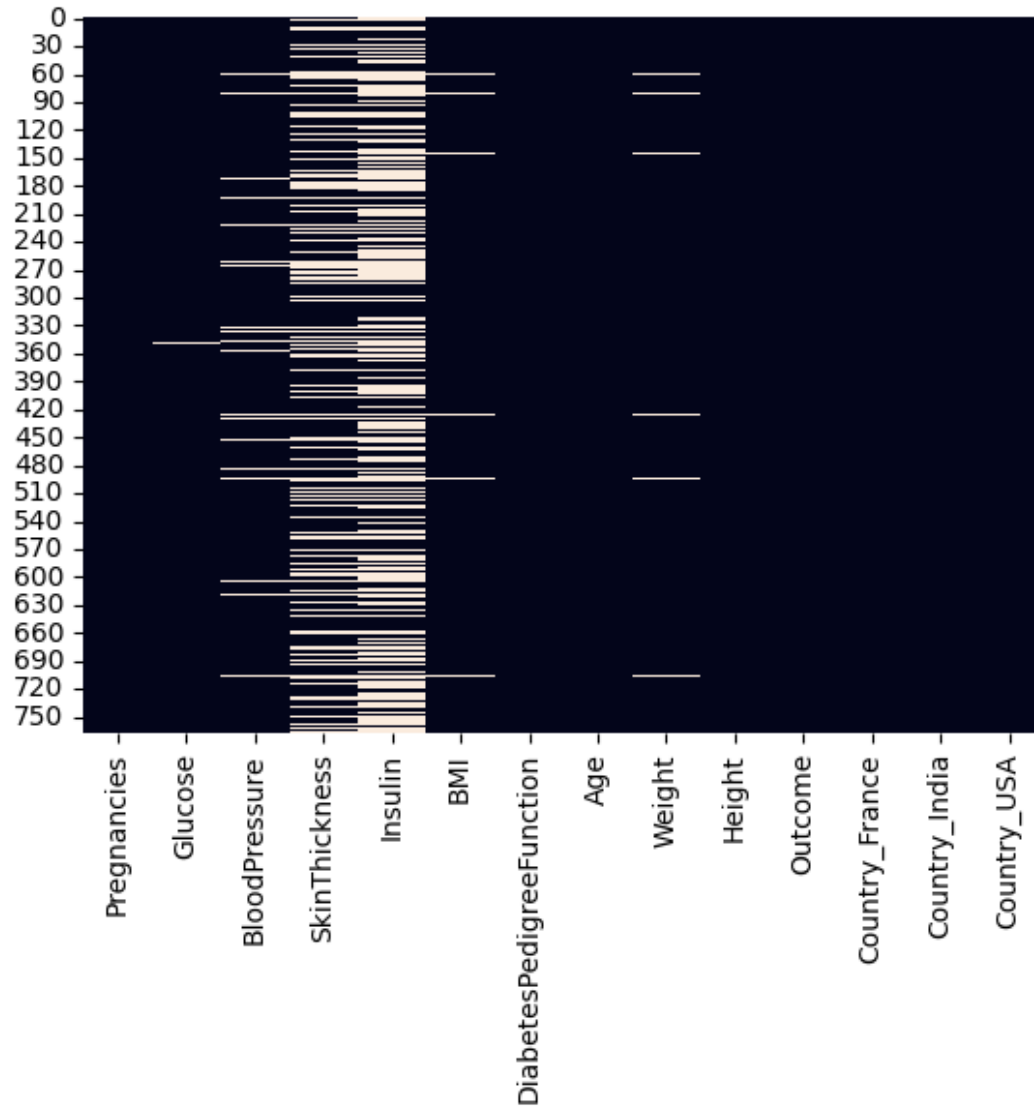
- **Niveau d'éducation** : secondaire, supérieur, doctorat
- **Niveau de risque** : Faible, Modéré, Elevé
- **Douleur ressentie** : Légère, Modérée, Sévère

## Label Encoding

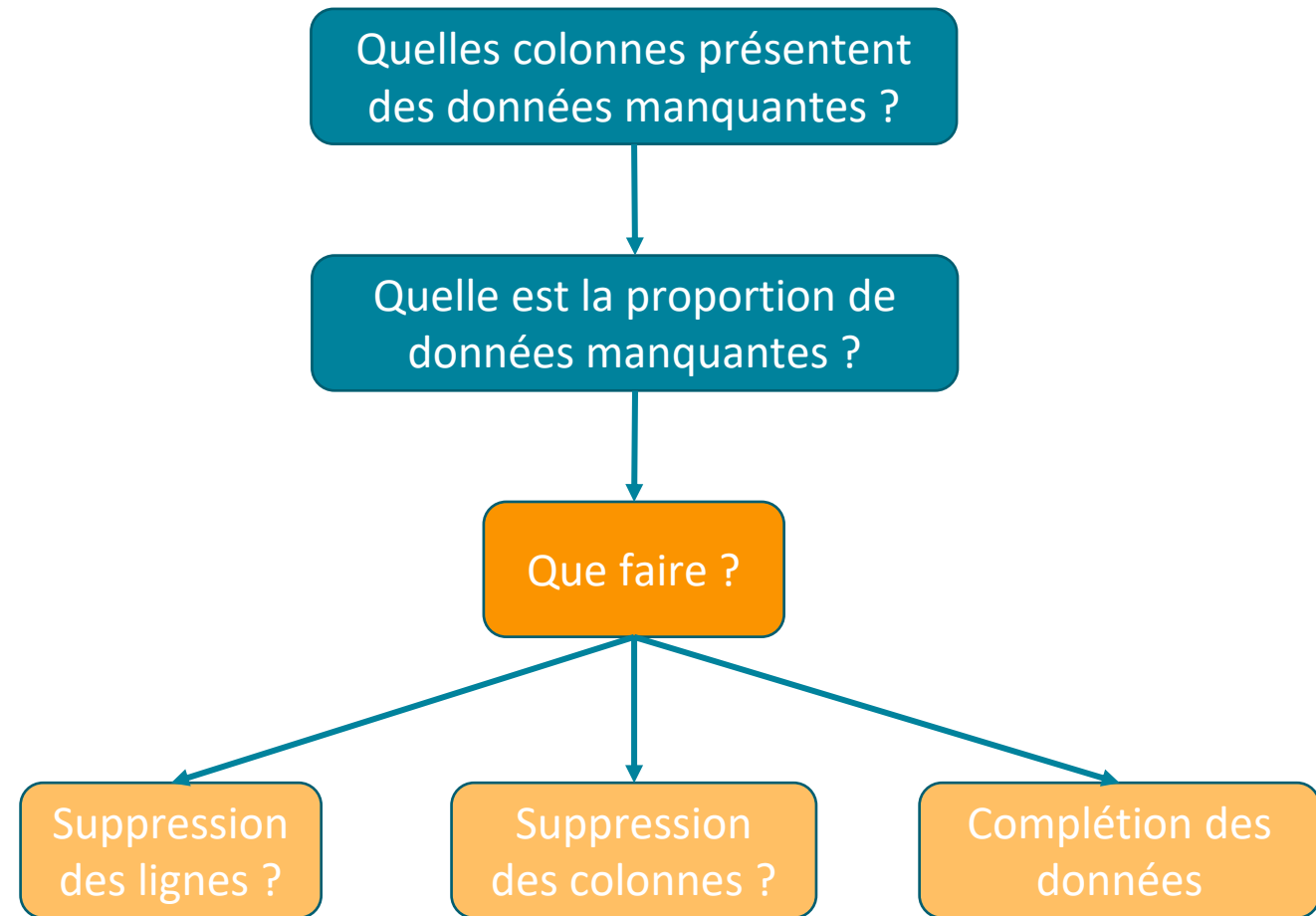
Douleur	Douleur
Légère	1
Sévère	3
Modérée	2
Sévère	3
NoPain	0
Légère	1

# Atelier 2 : Gestion des données manquantes

Complétude /  
Données manquantes



```
seaborn.heatmap(df_encoded.isna(), cbar=False)
```



# Atelier 2 : Gestion des données manquantes - Complétion

Complétude /  
Données manquantes

Remplacement par la  
moyenne

Feature B

10

X

25

$$X = (10 + 25) / 2$$

Interpolation linéaire  
(applicable sur des données  
ordonnées)

Feature A

8

11

12

Feature B

10

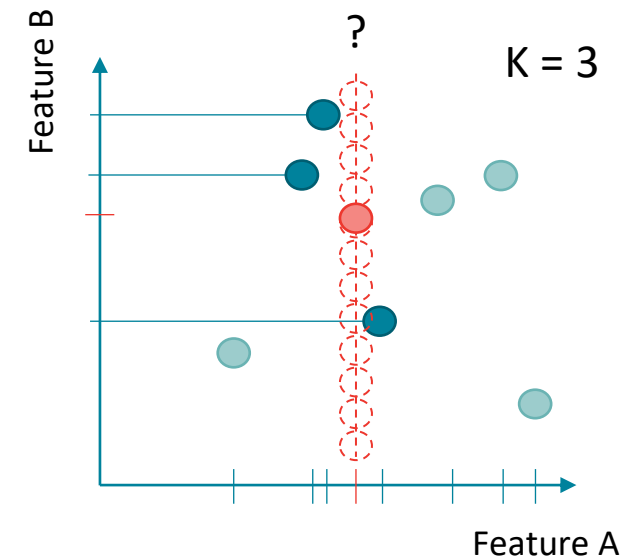
X

25

$$\text{pente} = (25 - 10) / (12 - 8)$$

$$X = 10 + \text{pente} * (11 - 8)$$

Complétion par plus  
proches voisins (K-NN)



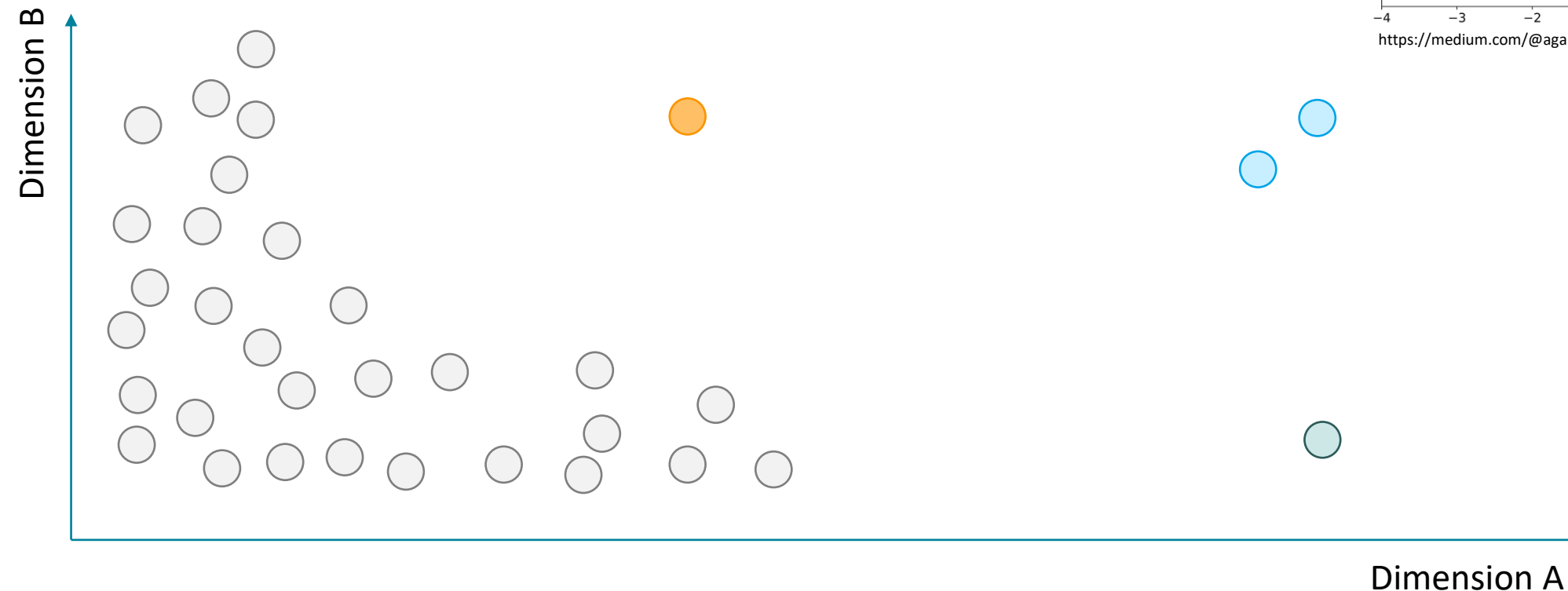
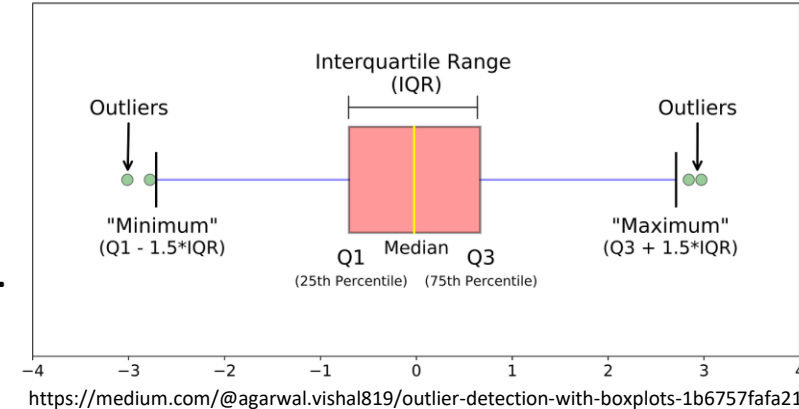
# Atelier 2 : Détection des outliers - Données aberrantes

**Outlier** (ou donnée aberrante) : Individu qui s'écarte significativement de la distribution initiale du jeu de données.

Certains algorithmes d'apprentissages sont très sensibles aux outliers.

On distingue :

- Les Outliers univariés, qui s'éloignent de la distribution sur une ou plusieurs dimensions de façon indépendante.
- Les Outliers multivariés, qui présentent une combinaison de features inhabituelle.



Source : <https://www.cuemath.com/data/outlier/>

# Atelier 2 : Équilibre de classes

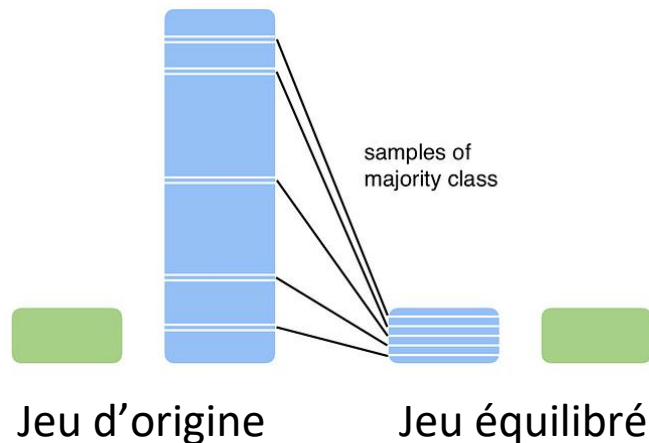
La performance d'un modèle est guidée par la question suivante :  
Le modèle a-t-il vu suffisamment d'exemples au cours de son entraînement ?

Hors, la majorité des jeux de données présentent un déséquilibre entre les différentes classes qu'il décrivent.

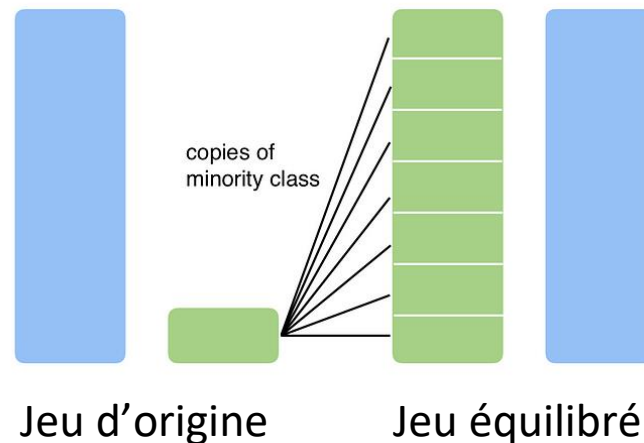
Un jeu de données idéal présente le même nombre d'individus pour chaque classe.



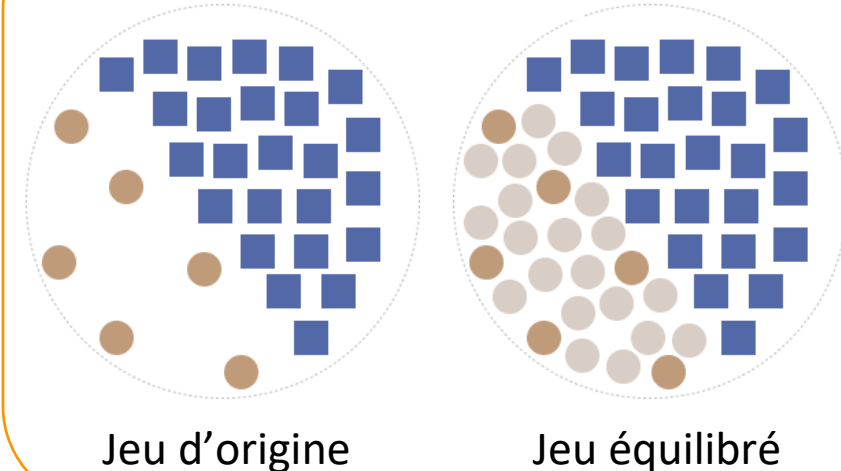
## Undersampling



## Oversampling



## SMOTE





Dans un jeu de données, il est possible d'avoir des features exprimées dans des échelles de données différentes

L'étape de normalisation permet de ramener les données sur une plage de valeurs comparable.

Cela permet de faciliter par la suite les calculs sur le jeu de données et de garantir une participation équitables des features à l'analyse et à l'apprentissage.

Pour normaliser les données on peut effectuer des transformations statistiques :

Normalisation Min-Max :

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Avec Python :

```
scaler = MinMaxScaler()  
scaler.fit_transform(features)
```

Standardisation :

$$x' = \frac{x - \mu}{\sigma}$$

Avec Python :

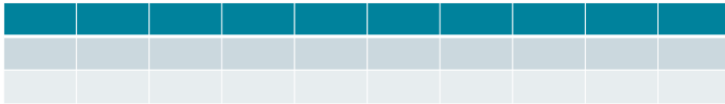
```
scaler = StandardScaler()  
scaler.fit_transform(features)
```



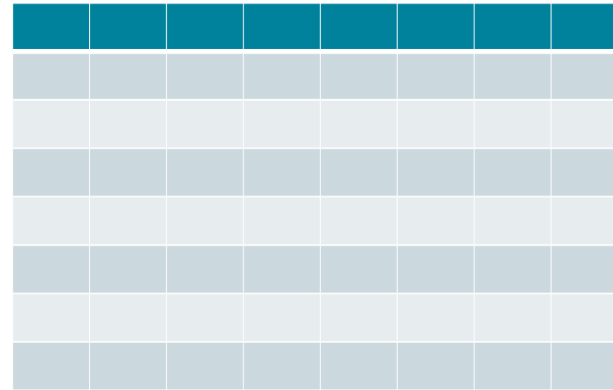
# Atelier 2 : Sélection des données

Pertinence et volume  
de données

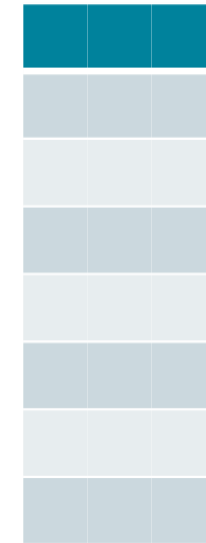
$N_{\text{individus}} \ll N_{\text{features}}$



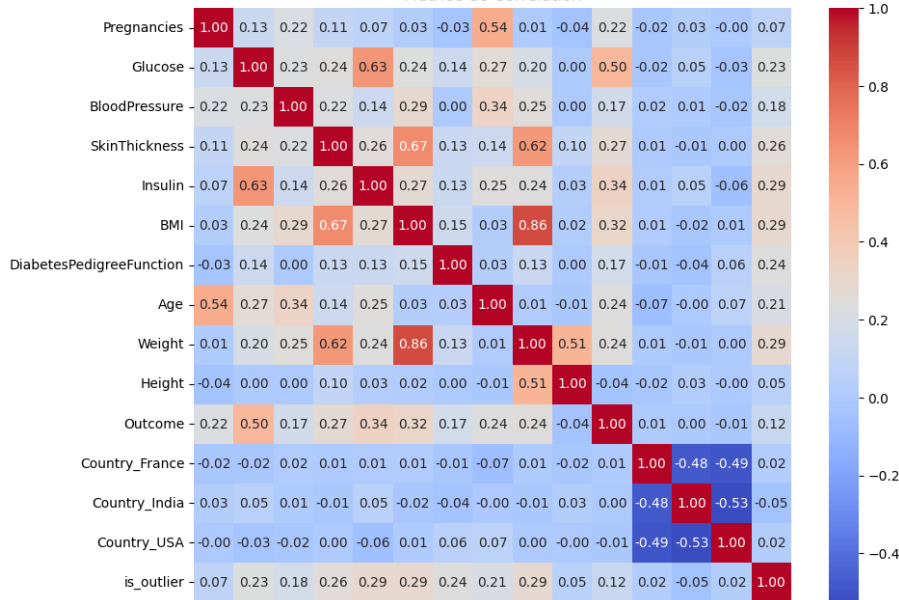
$N_{\text{individus}} = N_{\text{features}}$



$N_{\text{individus}} \gg N_{\text{features}}$



Matrice de Corrélation



Un outil pour la sélection  
de features : la matrice  
des corrélations

On cherche en général à  
avoir au moins 10  
individus pour 1 feature.